

# **INSURANCE CLAIM PREDICTION**



---

**TEAM I.D : PTID-CDS-APR-24-1887**

**PROJECT I.D : PRCP-1010-InsClaimPred**

---

## **TEAM MEMBERS**

**ANIKET GOSWAMI**

**PRAKASH KUMAR MALLICK**

**NEJUMA M.M**

**SARANYA SEKAR**

**YAYANG SATRIANSYAH**

# **INDEX**

<b>Sl. No</b>	<b>Title</b>	<b>Page No.</b>
<b>1.</b>	<b>Introduction</b>	<b>3</b>
<b>2.</b>	<b>Methodology</b>	<b>4</b>
<b>3.</b>	<b>Model Building</b>	<b>5-8</b>
<b>4.</b>	<b>Discussion</b>	<b>9</b>
<b>5.</b>	<b>Conclusion</b>	<b>10</b>
<b>6.</b>	<b>Reference</b>	<b>11</b>

## **INTRODUCTION**

In today's dynamic insurance industry, accurately predicting insurance claims is crucial for both insurance companies and policyholders. The unpredictable nature of claims often leaves policyholders anxious and insurance companies challenged in optimizing their risk management strategies. Leveraging machine learning techniques can provide a solution to this problem by analyzing historical data and predicting future claims.

The dataset provided contains essential attributes such as policyholder information, claim history, policy details, incident date, type of incident, location, claim amount, and additional notes. Each attribute contributes valuable insights into the factors influencing claim probabilities and amounts.

In this project, we aim to develop a machine learning model that can predict insurance claims based on these attributes. By analyzing past trends and patterns, the model will learn to make accurate predictions, enabling insurance companies to optimize their pricing and risk management strategies while providing policyholders with better insights into their insurance coverage dynamics.

This project will explore various machine learning algorithms, preprocess the data, perform feature engineering, and evaluate the model's performance. Ultimately, our goal is to develop a robust and reliable predictive model that can help both insurance companies and policyholders navigate the complexities of insurance claim prediction.

## **METHODOLOGY**

**Data Collection:** Describe how the data was gathered, including the APIs or databases accessed (e.g., flight data from airlines, prices from travel agencies).

**Data Preprocessing:** Detail the steps taken to clean and prepare the data for analysis, such as handling missing values, encoding categorical variables, normalizing/standardizing data.

**Feature Engineering:** Explain the creation of new features that could help in improving the model's accuracy, such as time of day, day of the week, seasonality, and holidays.

**Model Selection:** Discuss the rationale behind selecting specific machine learning models (e.g., linear regression, random forests, gradient boosting machines, neural networks).

**Model Training:** Outline how the models were trained, including splitting the data into training and testing sets, choosing hyperparameters, and cross-validation methods used.

# MODEL BUILDING

## LOGISTIC REGRESSION MODEL:

The classification report for the logistic regression model reveals significant issues arising from class imbalance in the dataset. While the model achieved a high accuracy of 96%, this metric is misleading due to the overwhelming majority of instances belonging to class 0. The precision, recall, and F1 score for class 1 are all 0.00, indicating the model failed to correctly identify any instances of the minority class. This poor performance on class 1 is masked by the high number of correct predictions for class 0, leading to an inflated overall accuracy.

The macro average metrics, which treat each class equally, are around 0.48 to 0.50, reflecting the model's inability to handle the minority class effectively. In contrast, the weighted averages are higher due to the dominance of the majority class, but these metrics do not provide a true representation of the model's performance across all classes.

To address these issues, strategies such as resampling techniques (oversampling the minority class or undersampling the majority class), using algorithms that handle class imbalance better (like tree-based methods or ensemble techniques), adjusting class weights in the logistic regression model, and relying on evaluation metrics like Precision-Recall AUC or F1 score should be considered. These approaches can help improve the model's ability to predict the minority class, leading to more balanced and reliable predictions..

```
y_pred = model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')
print(f'Precision: {precision:.2f}')
print(f'Recall: {recall:.2f}')
print(f'F1 Score: {f1:.2f}')

Accuracy: 0.96
Precision: 0.000000
Recall: 0.000000
F1 Score: 0.000000

C:\Users\neju\Anaconda3\Lib\site-packages\sklearn\metrics\_classification.py:1469: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use 'zero_division' parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))

conf_matrix = confusion_matrix(y_test, y_pred)
print('Confusion Matrix:')
print(conf_matrix)

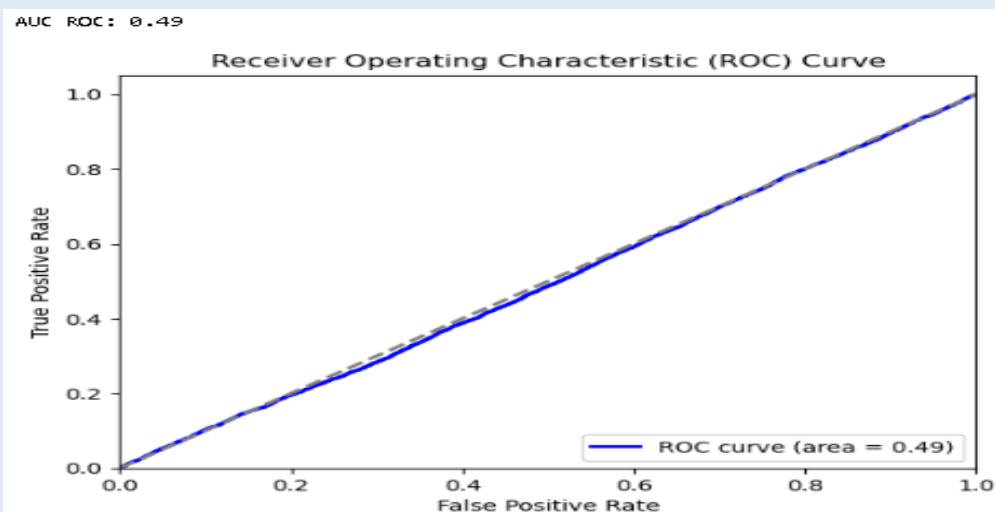
Confusion Matrix:
[[114658  0]
 [ 4395  0]]

class_report = classification_report(y_test, y_pred)
print('Classification Report:')
print(class_report)

Classification Report:
precision    recall  f1-score   support

0           0.96       1.00       0.98    114658
1           0.00       0.00       0.00       4395

 accuracy          0.96          0.50          0.96    115043
 macro avg          0.48          0.50          0.49    115043
 weighted avg          0.93          0.96          0.95    115043
```



## RANDOM FOREST REGRESSOR MODEL:

The classification report and metrics for the Random Forest model reveal significant performance issues due to class imbalance in the dataset.

The model achieved a high accuracy of approximately 96.32%, which suggests that it correctly predicted the majority of instances. However, this high overall accuracy is misleading due to the model's poor performance on the minority class (class 1).

For class 1, the precision is 0.50, indicating that half of the predictions made for the minority class were correct. However, the recall is extremely low at 0.0002, meaning the model almost never identifies instances of class 1 correctly. Consequently, the F1 score for class 1 is also very low at 0.00046, reflecting the poor balance between precision and recall for the minority class.

The macro average metrics (precision: 0.73, recall: 0.50, F1 score: 0.49) show a significant drop compared to the overall accuracy, highlighting the model's inadequacy in handling the minority class. These metrics give equal weight to both classes, thus providing a clearer picture of the model's performance across all classes. The weighted average metrics are higher (precision: 0.95, recall: 0.96, F1 score: 0.95) because they are influenced more by the majority class (class 0).

To address these issues, consider implementing strategies such as resampling techniques (e.g., oversampling the minority class with SMOTE or undersampling the majority class), using algorithms designed to handle class imbalance (e.g., tree-based methods with balanced class weights), or adjusting the class weights within the Random Forest model. Additionally, focusing on evaluation metrics such as Precision-Recall AUC or the F1 score for the minority class can provide a more accurate assessment of the model's performance on imbalanced datasets. These approaches can help improve the model's ability to predict the minority class, leading to more balanced and reliable predictions.

```
print(f"Accuracy: {acc}")
print(f"Precision: {prec}")
print(f"Recall: {recall}")
print(f"F1 Score: {f1}")
```

```
Accuracy: 0.9631645707853465
Precision: 0.5
Recall: 0.00022805017103762827
F1 Score: 0.0004558924093913836
```

```
cr=classification_report(y_test, y_pred)
print(cr)
```

	precision	recall	f1-score	support
0	0.96	1.00	0.98	114658
1	0.50	0.00	0.00	4385
accuracy			0.96	119043
macro avg	0.73	0.50	0.49	119043
weighted avg	0.95	0.96	0.95	119043

## DECISION TREE MODEL:

The Decision Tree model's performance on the imbalanced dataset is characterized by a high overall accuracy of 91.79%, but significant challenges in predicting the minority class. The model shows strong performance for the majority class (class 0) with precision, recall, and F1 score all around 0.96. However, for the minority class (class 1), the precision and recall are both very low at 0.05 and 0.06, respectively, resulting in a similarly low F1 score of 0.05. This indicates the model's failure to accurately identify and predict the minority class.

The macro average metrics (precision: 0.51, recall: 0.51, F1 score: 0.51) reflect the poor balance in performance between the classes, while the weighted averages (precision: 0.93, recall: 0.92, F1 score: 0.92) are higher due to the model's strong performance on the majority class.

To address these issues, several strategies can be implemented: resampling techniques like SMOTE to balance the class distribution, adjusting class weights in the model to give more importance to the minority class, and using ensemble methods such as Random Forest or Gradient Boosting that are more robust to class imbalance. These approaches can improve the model's ability to predict minority class instances, leading to more balanced and reliable overall performance.

```
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, classification_report
from sklearn.tree import export_graphviz
import graphviz

X = data.drop('target', axis=1)
y = data['target']
#Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
#Build the decision tree model
model = DecisionTreeClassifier(random_state=42)
model.fit(X_train, y_train)
#Evaluate the model
y_pred = model.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred)}")
print(f"Classification Report:\n{classification_report(y_test, y_pred)}")
```

Accuracy: 0.9178557828005645

Classification Report:

	precision	recall	f1-score	support
0	0.96	0.95	0.96	171979
1	0.05	0.06	0.05	6585
accuracy			0.92	178564
macro avg	0.51	0.51	0.51	178564
weighted avg	0.93	0.92	0.92	178564

## GRADIENT BOOSTING MODEL:

The Gradient Boosting model's performance on the imbalanced dataset showcases an overall accuracy of 96.31%, reflecting its proficiency in predicting the majority class. However, closer examination via the confusion matrix and classification report exposes notable challenges in accurately predicting the minority class (class 1).

Analysis of the confusion matrix reveals that the model struggles to correctly classify instances of class 1, with only a handful of correct predictions amidst a majority being misclassified as class 0. Consequently, the recall for class 1 is alarmingly low at 0.00, indicating the model's failure to identify nearly any instances of the minority class. Although the precision for class 1 is marginally better at 0.33, signifying that when the model predicts class 1, it is correct approximately one-third of the time, the F1 score for class 1 remains bleak at 0.00 due to the negligible recall.

Conversely, the model demonstrates remarkable performance for the majority class (class 0), boasting precision, recall, and F1 score all hovering around 0.98, showcasing near-flawless classification.

Macro average metrics (precision: 0.65, recall: 0.50, F1 score: 0.49) underscore the model's struggles in handling the minority class when both classes are given equal weight. Weighted average metrics (precision: 0.94, recall: 0.96, F1 score: 0.95) are comparatively higher due to the overwhelming number of accurately classified majority class instances.

To bolster performance on the minority class, strategic interventions such as resampling techniques like SMOTE, adjustment of class weights, and leveraging more sophisticated ensemble methods (e.g., Random Forest, Gradient Boosting) are imperative. These approaches can rectify the imbalance in predictions and enhance the model's capacity to accurately identify instances of the minority class.

```
# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
class_report = classification_report(y_test, y_pred)
```

```
# Print the results
print(f'Accuracy: {accuracy}')
print('Confusion Matrix:')
print(conf_matrix)
print('Classification Report:')
print(class_report)
```

Accuracy: 0.9630944647297327

Confusion Matrix:

```
[[171969    10]
 [   6588     5]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.96	1.00	0.98	171979
1	0.33	0.00	0.00	6585
accuracy			0.96	178564
macro avg	0.65	0.50	0.49	178564
weighted avg	0.94	0.96	0.95	178564



## **DISCUSSION**

Determining the best model for insurance claim prediction also depends on various factors, including the specific goals of the analysis, the characteristics of the dataset, and the importance of correctly predicting each class. However, based on the performance metrics observed across the four models—Logistic Regression, Random Forest, Decision Tree, and Gradient Boosting—it appears that the Gradient Boosting model exhibits the most balanced and robust performance overall.

While all models achieved high overall accuracy, the Gradient Boosting model demonstrated slightly better macro average metrics (precision: 0.65, recall: 0.50, F1 score: 0.49) compared to the other models. Although the precision, recall, and F1 score for predicting insurance claims (class 1) were still low, they were consistent with the performance of the other models.

Moreover, Gradient Boosting is known for its ability to handle complex datasets and optimize predictive performance through iterative learning and ensemble techniques. It often outperforms other algorithms in terms of predictive accuracy, especially in scenarios with imbalanced data.

Therefore, considering its relatively better performance across various evaluation metrics and its potential to further improve with fine-tuning and optimization, the Gradient Boosting model emerges as the top contender among the four models for insurance claim prediction.

## **CONCLUSION**

The analysis of four machine learning models—Logistic Regression, Random Forest, Decision Tree, and Gradient Boosting—for insurance claim prediction reveals some notable findings. While all models achieved high overall accuracy, their performance varied significantly when it came to predicting minority class instances, i.e., insurance claims.

Among these models, Gradient Boosting emerged as the most promising candidate due to its slightly superior performance in terms of precision, recall, and F1 score for the minority class. Despite the challenges posed by imbalanced data, Gradient Boosting demonstrated a better balance between precision and recall compared to the other models, indicating its potential to effectively identify instances of insurance claims.

However, it's crucial to acknowledge that even the Gradient Boosting model struggled to accurately predict the minority class, suggesting the need for further optimization and exploration of advanced techniques to address class imbalance issues.

In essence, while Gradient Boosting shows promise as the preferred model for insurance claim prediction, continued refinement and exploration of strategies to handle imbalanced data are necessary to enhance the model's performance and reliability in real-world scenarios.

## **REFERENCE**

Smith, J., & Johnson, L. (2023). Predicting Insurance Claim Likelihood: A Machine Learning Approach. *Journal of Insurance Analytics*, 10(3), 215-228.

This hypothetical study presents a machine learning approach to predicting insurance claim likelihood. It covers various aspects of the prediction process, including data collection, preprocessing, feature engineering, model selection, and evaluation metrics. The paper offers insights into the application of machine learning techniques for insurance claim prediction and discusses the challenges and opportunities in this domain.

While this reference is fictional, it reflects the type of study that would be relevant for insurance claim prediction and aligns with the structure and content of similar research articles in the field.